

# Traitement automatique du bambara

## Objectifs et premiers résultats

Damien Nouvel  
Valentin Vydrin  
Davy Auffret

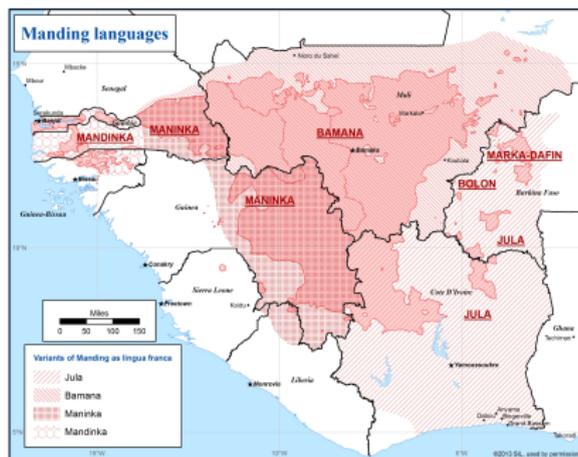
**ERTIM**  
Equipe de recherche  
textes, informatique,  
multilinguisme



# Plan

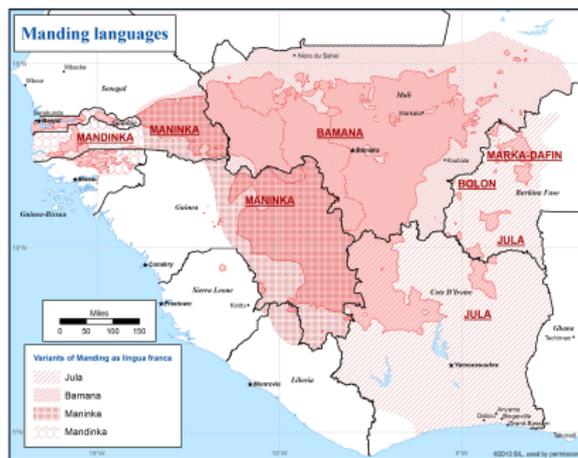
1. Le bambara
2. Projet MANTAL
3. Annotation morpho-syntaxique
4. Tonalisation du bambara
5. Conclusions et perspectives

# Contexte



- ▶ Parlée principalement au Mali (diglossie : français, 20%)
- ▶ Ou « bamanaka »

# Contexte



- ▶ Parlée principalement au Mali (diglossie : français, 20%)
- ▶ Ou « bamanaka »
- ▶ 4M de locuteurs (10M en 2<sup>ème</sup> langue)
- ▶ Véhiculaire, tradition orale
- ▶ Macro-langue mandingue (avec dioula, malinké, etc.)

# Quelques considérations linguistiques

- Prononciation : 7 voyelles, 20 consonnes, 3 tons

# Quelques considérations linguistiques

- ▶ Prononciation : 7 voyelles, 20 consonnes, 3 tons
- ▶ Alphabet
  - Langues mandingues : alphabet n'ko (1950, Unicode 5.0, rtl)
  - Bambara : alphabet latin
  - Depuis 1982 :  $\epsilon$  (U+025B/U+03B5),  $\text{ɔ}$ ,  $\eta$ ,  $\text{ɲ}$

⇒ Langues peu standardisées

# Quelques considérations linguistiques

- ▶ Prononciation : 7 voyelles, 20 consonnes, 3 tons
- ▶ Alphabet
  - Langues mandingues : alphabet n'ko (1950, Unicode 5.0, rtl)
  - Bambara : alphabet latin
  - Depuis 1982 :  $\epsilon$  (U+025B/U+03B5),  $\mathfrak{c}$ ,  $\eta$ ,  $\mathfrak{n}$

⇒ Langues peu standardisées

- ▶ Grammaire
  - Type : S AUX O V X, tonale
  - Pas de genre grammatical
  - Pas de conjugaison (marques prédicatives AUX)
  - Peu de flexion (-w : pluriel)

# Travaux antérieurs

## ▸ Références sur le bambara

- (Binger, 1886) Essai sur la langue bambara.
- (Sauvant, 1926) Dictionnaire bambara-français.
- (Vydrin, 1999) Les parties du discours en bambara.
- (Vydrin, 1999) Manding-english dictionary (maninka, bamana).
- (Dumestre, 2003) Grammaire fondamentale du bambara.
- (Bailleul, 2007) Dictionnaire bambara-français.
- (Vydrin, 2008) Glossed electronic corpora of Mande languages.
- (Dumestre, 2011) Dictionnaire bambara-français.
- (Vydrin, 2011) Corpus électronique annoté des textes bambara.
- (Enguehard, 2012) Vers l'info. de langues d'Afrique de l'Ouest.
- (Maslinsky, 2014) Daba : a model and tools for manding corpora.

# Plan

1. Le bambara
2. **Projet MANTAL**
3. Annotation morpho-syntaxique
4. Tonalisation du bambara
5. Conclusions et perspectives

# Corpus bambara de référence

- Collecte de textes en bambara
  - Publiés (périodiques, littérature) ou non (lettres, trans.)
  - Normalisation des textes (orthographe, tons, etc.)?
  - Textes en ligne <http://cormand.huma-num.fr/biblio/>

⇒ Volume : 2,3M mots

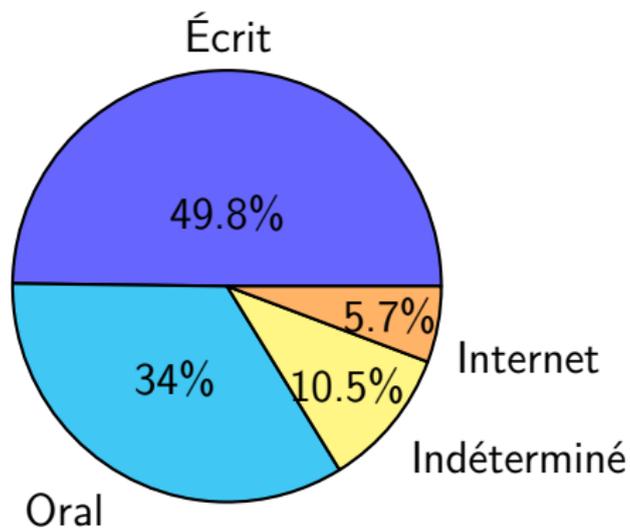
- Utilisation essentiellement linguistique
  - Apprentissage de la langue
  - Études linguistiques sur corpus
  - Annotation (parties du discours, lemmes, morphologie, gloses)

⇒ Labex EFL, axe 6 : ressources linguistiques

⇒ Site internet : <http://cormand.huma-num.fr/>

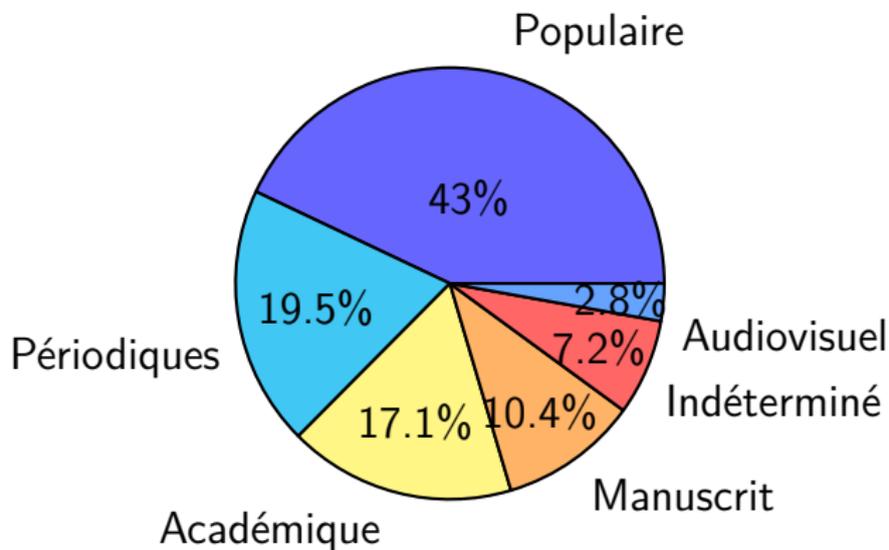
⇒ Modélisation linguistique informatisée (SketchEngine)

# Mediums du corpus



⇒ Prédominance de l'écrit, mais souvent issu de l'oral (contes, etc.)

# Sources du corpus



# Annotation du corpus

- Utilisation de **Daba** (Maslinsky, 2014)
  - Tokenisation
  - Recherche dans les dictionnaires
  - Analyses morphologiques

# Annotation du corpus

- Utilisation de **Daba** (Maslinsky, 2014)
  - Tokenisation
  - Recherche dans les dictionnaires
  - Analyses morphologiques

⇒ Pré-annotation automatique et ambiguë

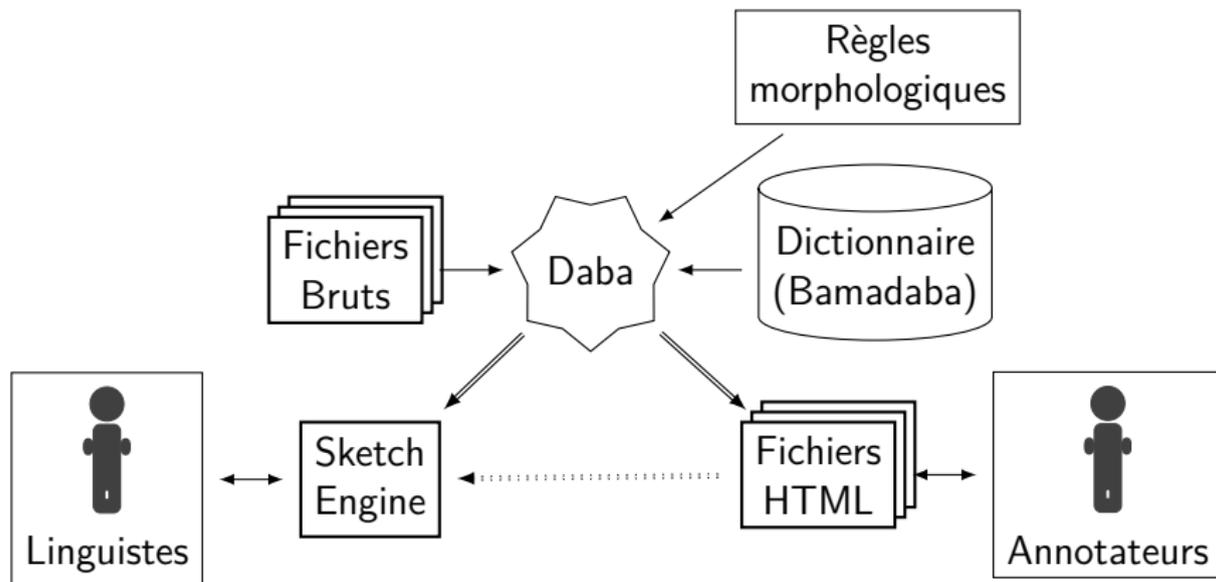
# Annotation du corpus

- Utilisation de **Daba** (Maslinsky, 2014)
  - Tokenisation
  - Recherche dans les dictionnaires
  - Analyses morphologiques

⇒ Pré-annotation automatique et ambiguë

- Annotation manuelle
  - Peu de moyens (profs, étudiants, bénévoles)
  - Validation humaine
  - Correction / normalisation

# Schéma de fonctionnement



# Objectifs de MANTAL

- Traitement des langues MANDingues avec des outils TAL
    - Normalisation / uniformisation
    - Parties du discours
    - (Entités nommées)
    - (Malinké)
- ⇒ Exploitation des données annotées (supervision)

# Objectifs de MANTAL

- ▶ Traitement des langues MANDingues avec des outils TAL
    - Normalisation / uniformisation
    - Parties du discours
    - (Entités nommées)
    - (Malinké)
- ⇒ Exploitation des données annotées (supervision)
- ▶ Cadre général
    - Collaboration LLACAN / ERTIM
    - Durée : 3 ans
    - Budget interne Inalco (stages, prestations, missions)
- ⇒ Interactions avec les linguistes

# Plan

1. Le bambara
2. Projet MANTAL
3. Annotation morpho-syntaxique
4. Tonalisation du bambara
5. Conclusions et perspectives

# Procédures d'annotation

- Plusieurs niveaux de traitement

# Procédures d'annotation

- Plusieurs niveaux de traitement
  - Fichier source

# Procédures d'annotation

- Plusieurs niveaux de traitement
  - Fichier source
  - Recherche des ponctuations, nombres et noms propres

# Procédures d'annotation

- Plusieurs niveaux de traitement
    - Fichier source
    - Recherche des ponctuations, nombres et noms propres
    - Utilisation de Daba
      - Recherche dans les dictionnaires
      - Analyse morphologique
- ⇒ Automatique, sortie ambiguë

# Procédures d'annotation

- Plusieurs niveaux de traitement
  - Fichier source
  - Recherche des ponctuations, nombres et noms propres
  - Utilisation de Daba
    - Recherche dans les dictionnaires
    - Analyse morphologique
  - ⇒ Automatique, sortie ambiguë
  - Annotation par les linguistes
  - ⇒ Ambiguïtés résiduelles

# Procédures d'annotation

- Plusieurs niveaux de traitement
  - Fichier source
  - Recherche des ponctuations, nombres et noms propres
  - Utilisation de Daba
    - Recherche dans les dictionnaires
    - Analyse morphologique
    - ⇒ Automatique, sortie ambiguë
  - Annotation par les linguistes
  - ⇒ Ambiguïtés résiduelles
- ⇒ Plusieurs versions de fichiers à synchroniser

## Fichier source

<h>Dijɛ Yaalala</h>

Nsiirin, nsiirin. N y'a bila den dɔ le kan.

Den nin ye sira dali a face ɛ, k'a b'a ɛ ka taga dijɛ yaala ka dɔ fara a hakili kan. A face ye sira d'a ma a ka taga yaala. A tagara yaala kɔɛbɛ.

A ɔɔɔla ka taga ben sogosu dɔ ma. A yɔɔ bɛɛ tolira, a ko : "E ! Ala bɛ se." Sogo nin wulila ka kum'a ɛ k'a kan'a ɔ ko Ala bɛ se, k'Ala ka se b'a ɲɛɛ. "A tɛmɛna sogo nin na ka taga ɲɛɛ, ka kɔɔn saba ye. Fɔɔ jalen bɛ, ji foyi t'a la. A filanan ji to kɔɔ. A sabanan ji b'o kɔɔ. A tɛmɛn'o la ka taga se cɛkɔɔnin dɔ ma. Cɛ nin kɔɔla kɔɛbɛ. [...]"

## Fichier Daba

Nsiirin, nsiirin.

Nsiirin , nsiirin .

nsíirin      nsíirin

n              n

conte        conte

N y'a bila den dɔ le kan.

N   y'      a   bila   den   dɔ   le   kan .

ń   y'      à   bila   dén   dɔ   le   kàn

pers   pm      pers   v      n      dtm   prt   pp

1SG   PFV.TR   3SG   mettre   enfant   certain   FOC   sur

## Fichier vertical

Token	Lemma	PdD	Glose	Compos.	Original	Tonal
nsiiri	nsiiri	n	conte		NSIIRI	nsiiri
naaninan	naaninan	ORD adj	quatrième	naani	NAANINAN	náaninan
dinye	dunya jinye  dununye dinye  jyen diyen	n	monde		Dɪɲɛ	dɪɲɛ
yaalala	yaalala	AG.PRM n		yaala	Yaalala	yáalala
nsiirin	nsiiri nsiirin	n	conte		Nsiirin	nsíirin
,	,	c	,	,	,	,
nsiirin	nsiiri nsiirin	n	conte		nsiirin	nsíirin
.	.	c	.	.	.	.
n	n	pers	1SG		N	ń
y'	ye y'	pm	PFV.TR		y'	y'
a	a	pers	3SG		a	à
bila	bil' bla bila	v	mettre		bila	bila
den	den	n	enfant		den	dén
do	do	dtm	certain		dɔ	dó
le	le	prt	FOC		le	le
kan	kan	pp	sur		kan	kàn
.	.	c	.	.	.	.

# Volumétrie des ressources

- Corpus

<b>Corpus</b>	<b>Balises</b>	<b>Ponctuations</b>	<b>Formes (distinctes)</b>
Brut	412K	383K	2 321K (68K)
Désambiguïsé	104K	71K	426K (19K)

# Volumétrie des ressources

## ▸ Corpus

Corpus	Balises	Ponctuations	Formes (distinctes)
Brut	412K	383K	2 321K (68K)
Désambiguïsé	104K	71K	426K (19K)

## ▸ Dictionnaires (disponibles en ligne)

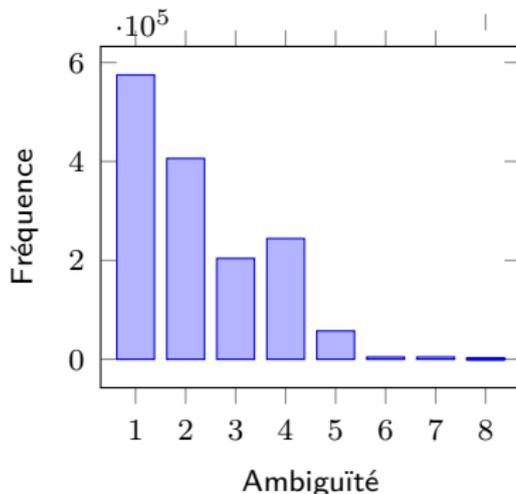
Dictionnaire	Description	Entrées	Ambiguïté
bamadaba	Dictionnaire principal	11K	1,137
enciclop	Notions encyclopédiques	29	1
jamuw	Noms claniques	375	1,001
togow	Prénoms	496	1
yorow	Toponymes	299	1

# Statistiques sur les ambiguïtés

- Corpus échantillon (1,5M mots)
  - Traité avec Daba
  - 1,285M mots (215K ponctuations)
  - 72K (5,6%) parties du discours non renseignées
  - 288K (22%) non-ambiguës (selon Daba)

# Statistiques sur les ambiguïtés

- ▶ Corpus échantillon (1,5M mots)
  - Traité avec Daba
  - 1,285M mots (215K ponctuations)
  - 72K (5,6%) parties du discours non renseignées
  - 288K (22%) non-ambiguës (selon Daba)



Ambiguïté	Fréquence
pers intj	75K
n v	72K
cop pp v pm	57K
adv conj pm v	53K
cop pm	43K
v n	38K
prt conj prn	33K
n.prop pp conj intj	23K
conj n prep v	23K
vq v adj n cop	22K

# Caractérisation des ambiguïtés

- Ambiguïtés fréquentes

Mot	Ambiguïté	Fréquence
à	pers intj	74K
yé	cop pp v pm	58K
kà	adv conj pm v	53K
ò	prt conj prn	34K

# Caractérisation des ambiguïtés

- Ambiguïtés fréquentes

Mot	Ambiguïté	Fréquence
à	pers intj	74K
yé	cop pp v pm	58K
kà	adv conj pm v	53K
ò	prt conj prn	34K

- Par mots distincts

Ambiguïté	Fréquence
PL n	1304
n v	1189
v n	1039
ptcp PTCP.RES	588
PFV.INTR v	584
NMLZ n	568
n.prop dtm ptcp prn adj PL n	497
n.prop n	483

# Jeu d'étiquettes sur le corpus désambiguïsé

Code	Partie du discours	Quantité
n	nom	82K
c	ponctuation	66K
pers	pronom personnel	54K
v	verbe	51K
pm	marque prédicative	41K
pp	postposition	34K
conj	conjonction	21K
cop	copule	18K
n.prop	nom propre	12K
dtm	déterminatif	12K
prn	pronom (non-personnel)	10K
prt	particule	10K

num	numératif	6K
adj	adjectif	4K
ptcp	participe	4K
intj	interjection	2K
adv	adverbe	2K
vq	verbe qualitatif	1K
onomat	onomatopée	102
adv.p	adverbe préverbal	26
conv.n	converbe nu	24
mrph	morphème	13

⇒ Jeu d'étiquettes relativement standard

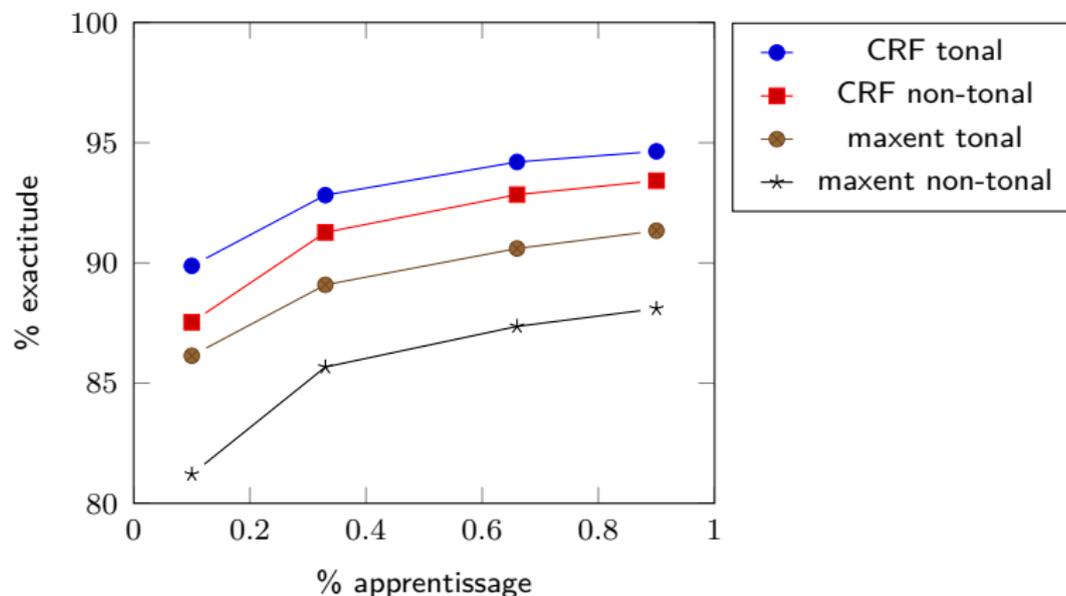
# Apprentissage : algorithmes et données

- Méthodologie
  - Utilisation de Wapiti (Lavergne, 2010)
  - Validation croisée à 10 plis

# Apprentissage : algorithme et données

## ▸ Méthodologie

- Utilisation de Wapiti (Lavergne, 2010)
- Validation croisée à 10 plis



# Médiums avec corpus complets

- Méthodologie
  - Sélection une sous-partie pour l'apprentissage
  - Évaluation une autre sous-partie (sans validation croisée)

# Médiums avec corpus complets

- Méthodologie
  - Sélection une sous-partie pour l'apprentissage
  - Évaluation une autre sous-partie (sans validation croisée)

App \ Test	Écrit	Internet	Oral	Indet.
Écrit	<i>98.55</i>	92.0	90.36	<b>92.93</b>
Internet	83.78	<i>99.08</i>	82.63	<b>84.20</b>
Oral	88.43	88.92	<i>98.81</i>	<b>90.68</b>
Indet.	87.26	<b>88.26</b>	87.07	<i>99.21</i>

# Médiums avec corpus équilibrés

- Méthodologie
  - Sélection une sous-partie pour l'apprentissage
  - Sous-échantillonnage selon corpus d'évaluation :  
25K éléments à la phrase près
  - Évaluation une autre sous-partie

# Médiums avec corpus équilibrés

## ▸ Méthodologie

- Sélection une sous-partie pour l'apprentissage
- Sous-échantillonnage selon corpus d'évaluation : 25K éléments à la phrase près
- Évaluation une autre sous-partie

App \ Test	Écrit	Internet	Oral	Indet.
Écrit	<i>98.97</i>	<b>86.09</b>	85.79	<b>86.09</b>
Internet	85.52	<i>99.08</i>	84.22	82.68
Oral	85.56	85.24	<i>99.08</i>	85.78
Indet.	<b>85.92</b>	84.70	<b>86.27</b>	<i>99.36</i>

# Traits pour l'apprentissage

- Ajouts de traits pour
  - Préfixe de 2 et 3 caractères
  - Suffixe de 2 et 3 caractères
  - Une version transformé du mot par un expression régulière, ANP
  - Une version non-tonalisée du mot
  - La taille du mot
  - Les étiquettes possibles dans les dictionnaires

# Traits pour l'apprentissage

- Ajouts de traits pour
  - Préfixe de 2 et 3 caractères
  - Suffixe de 2 et 3 caractères
  - Une version transformé du mot par un expression régulière, ANP
  - Une version non-tonalisée du mot
  - La taille du mot
  - Les étiquettes possibles dans les dictionnaires

Score	Base	Non-tonal	Préf.	Suf.	Dico	ANP	Taille	Tous
<b>Mot U</b>	86.14	86.36	89.09	89.46	89.87	89.41	89.14	90.65
<b>Phrase U</b>	18.70	19.09	24.38	25.60	25.62	25.18	23.82	29.04
<b>Mot B</b>	85.94	91.27	91.69	91.91	91.50	86.84	88.37	94.22
<b>Phrase B</b>	13.14	35.35	36.39	37.51	33.94	15.20	21.57	47.90

# Comparaison avec TreeTagger

- Tests TreeTagger
  - Entraînement : 90% du corpus, validation croisée
  - Pas de ressource additionnelle
  - Configuration par défaut

# Comparaison avec TreeTagger

- Tests TreeTagger
  - Entraînement : 90% du corpus, validation croisée
  - Pas de ressource additionnelle
  - Configuration par défaut

<b>Outil</b>	<b>Score</b>
Baseline	22%
Majorité	82,06%
TreeTagger	93,50
Wapiti	94.22
Majorité (non-tonal)	79,68%

# Plan

1. Le bambara
2. Projet MANTAL
3. Annotation morpho-syntaxique
4. Tonalisation du bambara
5. Conclusions et perspectives

# Utilisation des tons

- ▶ Caractéristiques

- Trois marques tonales : ` , ´ , ˇ (caron, hatchek)

⇒ Change le sens du mot

- Exemples :

- bá = maman / bà = chèvre
- tugu = bras / túgu = fermer
- tà = prendre, porter / tá = feu, propriété

- ▶ Les tons sont peu souvent marqués à l'écrit

- ▶ Essentiellement en 1<sup>ère</sup> syllabe

# Utilisation des tons

- ▶ Caractéristiques
  - Trois marques tonales : ` , ´ , ˇ (caron, hatchek)
  - ⇒ Change le sens du mot
  - Exemples :
    - bá = maman / bà = chèvre
    - tùgu = bras / túgu = fermer
    - tà = prendre, porter / tá = feu, propriété
- ▶ Les tons sont peu souvent marqués à l'écrit
- ▶ Essentiellement en 1<sup>ère</sup> syllabe
- ⇒ La présence de tons aide pour la morpho-syntaxe
- ⇒ Détecter automatiquement les tons ?

# Statistiques sur les tons

- Méthodologie
  - Sélection du corpus en version tonale (traitée)
  - Suppression des tons et comparaison des versions non-tonales
  - Tests de « re-tonalisation »

# Statistiques sur les tons

- Méthodologie

- Sélection du corpus en version tonale (traitée)
- Suppression des tons et comparaison des versions non-tonales
- Tests de « re-tonalisation »

⇒ 17 335 formes non-tonales, en moyenne 1.11 tonalisations

# Statistiques sur les tons

- Méthodologie

- Sélection du corpus en version tonale (traitée)
- Suppression des tons et comparaison des versions non-tonales
- Tests de « re-tonalisation »

⇒ 17 335 formes non-tonales, en moyenne 1.11 tonalisations

⇒ 1 518 tonalisations ambiguës, en moyenne 2,26 tonalisations

# Statistiques sur les tons

## ▸ Méthodologie

- Sélection du corpus en version tonale (traitée)
- Suppression des tons et comparaison des versions non-tonales
- Tests de « re-tonalisation »

⇒ 17 335 formes non-tonales, en moyenne 1.11 tonalisations

⇒ 1 518 tonalisations ambiguës, en moyenne 2,26 tonalisations

⇒ Baseline, tonalisation la plus fréquente : 59% (dont ponctuations)

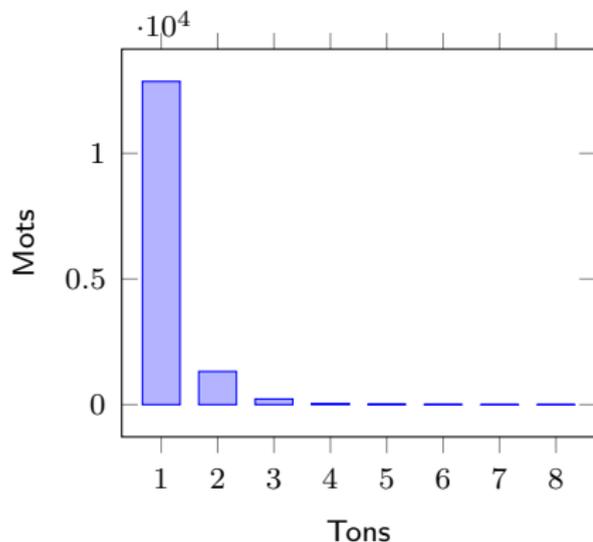
# Répartitions des tons

- Méthodologie
  - Répartition du nombre de tonalisations dans le corpus
  - Comparaison : mots / occurrences de mots

# Répartitions des tons

## ► Méthodologie

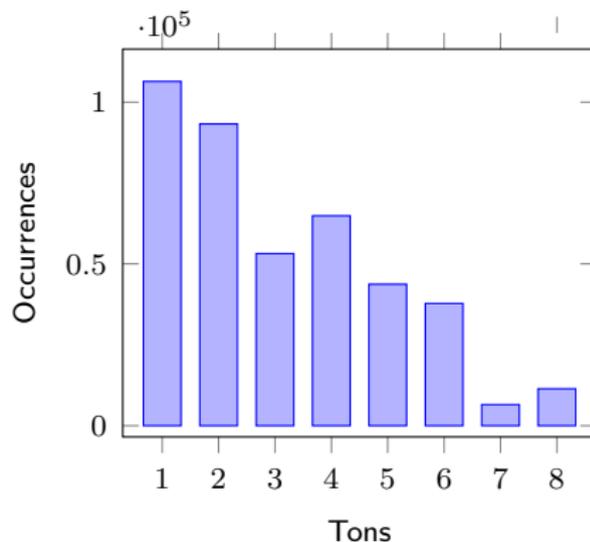
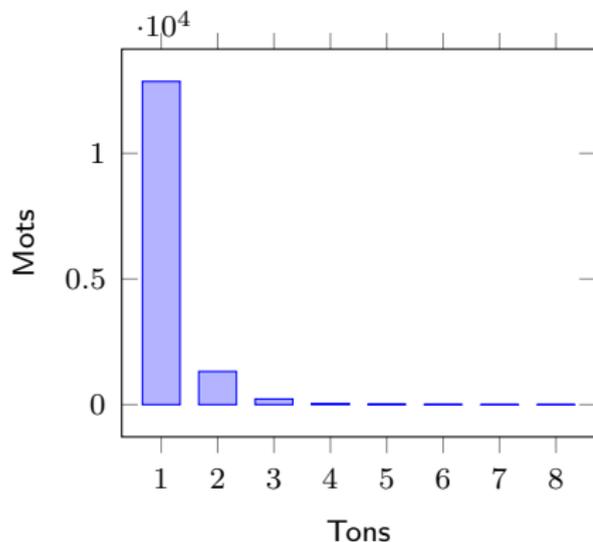
- Répartition du nombre de tonalisations dans le corpus
- Comparaison : mots / occurrences de mots



# Répartitions des tons

## ► Méthodologie

- Répartition du nombre de tonalisations dans le corpus
- Comparaison : mots / occurrences de mots



# Entropie des tons

- Méthodologie
  - Probabilités de tonalisations
  - Calcul d'entropie par mot

# Entropie des tons

- Méthodologie
  - Probabilités de tonalisations
  - Calcul d'entropie par mot

Ent.	Tonalisations	Traduction
3.20	táatúmà : 1.0 , tàatúmá : 1.0 , táatúmá : 1.0 , táatúma : 9.0 , táatúma : 1.0 , tàatúmà : 1.0 , tàatúma : 1.0 , tàatùmà : 1.0 , tàatùmá : 1.0 , táatùmá : 1.0 , táatùmà : 2.0 , tàatuma : 1.0 , tàatùma : 1.0 , táatuma : 1.0	Départ (?)
2.50	bámànanke : 10.0 , bàmànakè : 16.0 , bàmànanke : 5.0 , bàmànanke : 6.0 , cè : 16.0 , bàmànanke : 1.0 , bàmànan : 16.0 , bàmànanke : 1.0	Bambara
2.25	súurun : 1.0 , súuru : 1.0 , sùruntu : 2.0 , sùruntu : 1.0 , sùrundu : 1.0	Verser
2.02	cíyèn : 9.0 , tíjɛ : 11.0 , tíjɛ : 2.0 , ciyén : 2.0 , ciyèn : 11.0	Vérité
2.00	jènen : 1.0 , jénen : 1.0 , jànen : 1.0 , jè : 1.0	Regard
1.99	jógɔri : 7.0 , nwàri : 1.0 , jógɔri : 1.0 , jógɔri : 5.0 , júari : 5.0	Approcher
1.99	jènen : 1.0 , jénen : 1.0 , jànen : 1.0 , jè : ? 1.0	?
1.99	jógɔri : 7.0 , nwàri : 1.0 , jógɔri : 1.0 , jógɔri : 5.0 , júari : 5.0	Salir
1.99	lé : 1.0 , lè : 4.0 , lè : 1.0 , le : 7.0 , dè : 4.0	(clan)
1.95	tàamasyenw : 2.0 , tàamashyèn : 2.0 , táamashyèn : 1.0 , tàamashyèn : 2.0	Indiquer

# Plan

1. Le bambara
2. Projet MANTAL
3. Annotation morpho-syntaxique
4. Tonalisation du bambara
5. Conclusions et perspectives

# Conclusions

- ▶ Corpus du bambara
  - Initiative déjà ancienne
  - Corpus plutôt TAL-compatible
- ▶ Analyse en parties du discours
  - Assez bonnes performances
  - Problématique de tonalisation
- ▶ Tonalisation
  - Utile à la morpho-syntaxe
  - Probablement utile pour les lemmes et gloses

# Perspectives

- ▶ Travaux en cours
  - Coupler tonalisation et morpho-syntaxe
  - Essais de lemmatisation (jointe?)
  - Statistiques et expériences avec les gloses
- ▶ Entités nommées (2<sup>ème</sup> année)
  - Annotation en entités nommées
  - Translittération / transcription d'entités nommées
  - Expressions composées
- ▶ Malinké (3<sup>ème</sup> année)
  - Travail sur un autre corpus
  - Proximité avec le bambara